

CLIP2Scene: Towards Label-efficient 3D Scene Understanding by CLIP

Runnan Chen¹, Youquan Liu², Lingdong Kong³, Xinge Zhu⁶, Yuexin Ma⁵,
Yikang Li⁴, Yuenan Hou⁴, Yu Qiao⁴, Wenping Wang⁷

¹The University of Hong Kong

²Hochschule Bremerhaven

³National University of Singapore

⁴Shanghai AI Lab

⁵ShanghaiTech University

⁶The Chinese University of Hong Kong

⁷Texas A&M University

Abstract

Contrastive language-image pre-training (CLIP) achieves promising results in 2D zero-shot and few-shot learning. Despite the impressive performance in 2D tasks, applying CLIP to help the learning in 3D scene understanding has yet to be explored. In this paper, we make the first attempt to investigate how CLIP knowledge benefits 3D scene understanding. To this end, we propose CLIP2Scene, a simple yet effective framework that transfers CLIP knowledge from 2D image-text pre-trained models to a 3D point cloud network. We show that the pre-trained 3D network yields impressive performance on various downstream tasks, i.e., annotation-free and fine-tuning with labelled data for semantic segmentation. Specifically, built upon CLIP, we design a Semantic-driven Cross-modal Contrastive Learning framework that pre-trains a 3D network via semantic and spatial-temporal consistency regularization. For semantic consistency regularization, we first leverage CLIP’s text semantics to select the positive and negative point samples and then employ the contrastive loss to train the 3D network. In terms of spatial-temporal consistency regularization, we force the consistency between the temporally coherent point cloud features and their corresponding image features. We conduct experiments on the nuScenes and SemanticKITTI datasets. For the first time, our pre-trained network achieves annotation-free 3D semantic segmentation with 20.8% mIoU. When fine-tuned with 1% or 100% labelled data, our method significantly outperforms other self-supervised methods, with improvements of 8% and 1% mIoU, respectively. Furthermore, we demonstrate its generalization capability

Semantic-driven Cross-modal Contrastive Learning

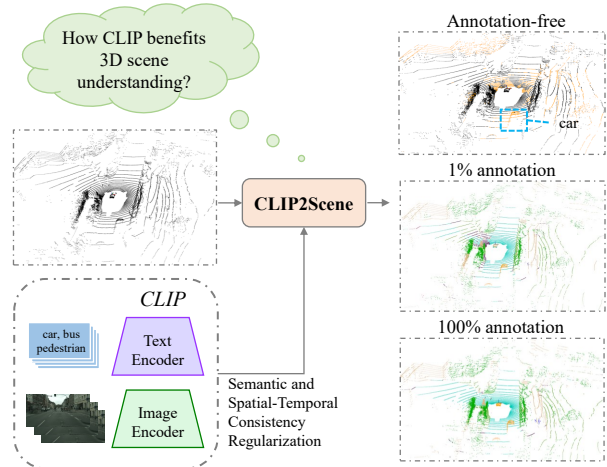


Figure 1. We explore whether and how CLIP knowledge benefits 3D scene understanding. To this end, we propose CLIP2Scene, a semantic-driven cross-modal contrastive learning framework that leverages CLIP knowledge to pre-train a 3D point cloud segmentation network via semantic and spatial-temporal consistency regularization. CLIP2Scene yields impressive performance on annotation-free 3D semantic segmentation and significantly outperforms other self-supervised methods when fine-tuning on annotated data.

for handling cross-domain datasets.

1. Introduction

3D scene understanding is fundamental in autonomous driving, robot navigation, etc [24, 26]. Current deep

learning-based methods have shown inspirational performance on 3D point cloud data [37, 50, 29, 44, 15, 45]. However, some drawbacks hinder their real-world applications. The first one comes from their heavy reliance on the large collection of the annotated point clouds, especially when high-quality 3D annotations are expensive to acquire [34, 40]. Besides, they typically fail to recognize novel objects that are never seen in the training data [11, 35]. As a result, it may need extra annotation efforts to train the model on recognizing these novel objects, which is both tedious and time-consuming.

Contrastive Vision-Language Pre-training (CLIP) [38] provides a new perspective that mitigates the above issues in 2D vision. It was trained on large-scale free-available image-text pairs from websites and built vision-language correlation to achieve promising open-vocabulary recognition. MaskCLIP [49] further explores semantic segmentation based on CLIP. With minimal modifications to the CLIP pre-trained network, MaskCLIP can be directly used for the semantic segmentation of novel objects without additional training efforts. PointCLIP [48] reveals that the zero-shot classification ability of CLIP can be generalized from the 2D image to the 3D point cloud. It respectively projects a point cloud frame into different views of 2D depth maps that bridge the modal gap between the image and the point cloud. The above studies indicate the potential of CLIP on enhancing the 2D segmentation and 3D classification performance. However, whether and how CLIP knowledge benefits 3D scene understanding is still under-explored.

In this paper, we explore how to leverage CLIP’s 2D image-text pre-learned knowledge for 3D scene understanding. Previous cross-modal knowledge distillation methods [40, 34] suffer from the optimization-conflict issue, *i.e.*, some of the positive pairs are regarded as negative samples for contrastive learning, leading to unsatisfactory representation learning and hammering the performance of downstream tasks. Besides, they also ignore the temporal coherence of the multi-sweep point cloud, failing to utilize the rich inter-sweep correspondence. To handle the mentioned problems, we propose a novel Semantic-driven Cross-modal Contrastive Learning framework that fully leverages CLIP’s semantic and visual information to regularize a 3D network. Specifically, we propose Semantic Consistency Regularization and Spatial-Temporal Consistency Regularization. In semantic consistency regularization, we utilize CLIP’s text semantics to select the positive and negative point samples for less-conflict contrastive learning. For spatial-temporal consistency regularization, we take CLIP’s image pixel feature to impose a soft constraint on points within local space and time. Such operation also prevents the network from degenerating due to image-to-point calibration errors.

We conduct several downstream tasks on nuScenes to verify how the pre-trained network benefits the 3D scene understanding. The first one is annotation-free semantic segmentation. Following MaskCLIP, we place class names into multiple hand-crafted templates as prompts and average the text embeddings generated by CLIP to conduct the annotation-free segmentation. For the first time, our method achieves 20.8% mIoU annotation-free 3D semantic segmentation without any labelled data for training. Secondly, we compare with other self-supervised methods to verify the superiority of our method in label-efficient learning. When fine-tuning the 3D network with 1% or 100% labelled data, our method significantly outperforms state-of-the-art self-supervised methods, with improvements of 8% and 1% mIoU, respectively. Besides, to verify the generalization capability, we pre-train the network on the nuScenes dataset and evaluate it on the SemanticKITTI dataset. Our method still significantly outperforms state-of-the-art methods.

The contributions of our work are summarized as follows.

- The first work that distills CLIP knowledge to a 3D network for 3D scene understanding.
- We propose a novel Semantic-driven Cross-modal Contrastive Learning framework that pre-trains a 3D network via spatial-temporal and semantic consistency regularization.
- We propose a novel Semantic-guided Spatial-Temporal Consistency Regularization that forces the consistency between the temporally coherent point cloud features and their corresponding image features.
- For the first time, our method achieves promising performance on annotation-free 3D scene segmentation and significantly outperforms state-of-the-art self-supervised methods when fine-tuning with labelled data.

2. Related Work

Zero-shot Learning in 3D. The objective of zero-shot learning (ZSL) is to recognize objects that are unseen in the training set. Many efforts have been devoted to the 2D recognition tasks [8, 30, 47, 36, 31, 1, 43, 32, 4, 2, 19, 33, 23], and few works concentrate on performing ZSL in the 3D domain [18, 11, 35, 16, 17]. [18] makes the first attempt to apply ZSL to 3D tasks, where they train PointNet [37] on “seen” samples and test on “unseen” samples. Subsequent work [16] addresses the hubness problem caused by the low-quality point cloud features. [17] proposes the triplet loss to boost the performance under the transductive setting, where the “unseen” class is observed and unlabeled

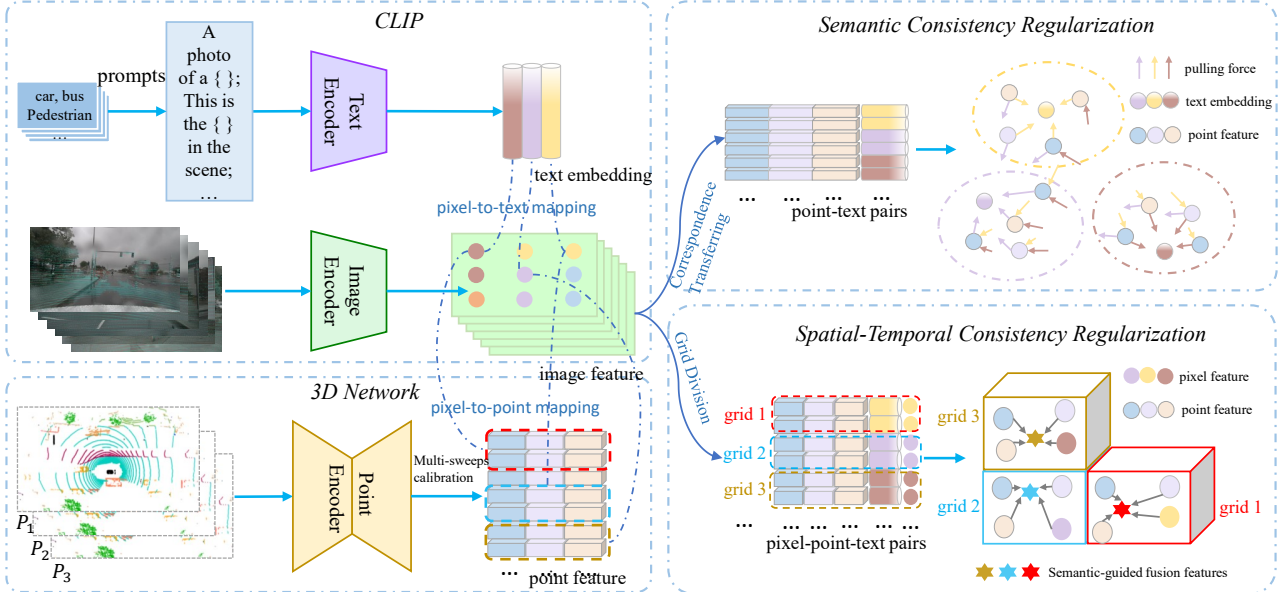


Figure 2. Illustration of the Semantic-driven Cross-modal Contrastive Learning. Firstly, we obtain the text embeddings t_i , image pixel feature x_i , and point feature p_i by text encoder, image encoder, and point encoder, respectively. Secondly, we leverage CLIP knowledge to construct positive and negative samples for contrastive learning. Thus we obtain point-text pairs $\{x_i, t_i\}_{i=1}^M$ and all pixel-point-text pairs in a short temporal $\{\hat{x}_i^k, \hat{p}_i^k, t_i^k\}_{i=1, k=1}^{M, K}$. Here, $\{x_i, t_i\}_{i=1}^M$ and $\{\hat{x}_i^k, \hat{p}_i^k, t_i^k\}_{i=1, k=1}^{M, K}$ are used for Semantic Consistency Regularization and Spatial-Temporal Consistency Regularization, respectively. Lastly, we perform Semantic Consistency Regularization by pulling the point features to their corresponding text embedding and Spatial-Temporal Consistency Regularization by mimicking the temporally coherent point features to their corresponding pixel features.

in the training phase. Recently, some studies introduced CLIP into zero-shot learning. MaskCLIP [49] investigates the problem of utilizing CLIP to help the 2D dense prediction tasks and exhibits encouraging zero-shot semantic segmentation performance. PointCLIP [48] is the pioneering work that applies CLIP to 3D recognition. As opposed to previous approaches that require training on the labelled point cloud, PointCLIP is free from any 3D training and shows impressive performance on zero-shot and few-shot classification tasks. Our work takes a step further to investigate whether the rich semantic and visual knowledge in CLIP can benefit the 3D semantic segmentation tasks.

Self-supervised Representation Learning. The purpose of self-supervised learning is to learn a good representation that benefits the downstream tasks. The dominant approaches resort to contrastive learning to pre-train the network [27, 25, 21, 20, 14, 13, 7, 10, 12, 9]. Recently, inspired by the success of CLIP, leveraging the pre-trained model of CLIP to the downstream tasks has raised the community’s attention. DenseCLIP [39] utilizes the CLIP’s pre-trained knowledge for dense image pixel prediction. DetCLIP [46] proposes a pre-training method equipped with CLIP for open-world detection. In this paper, we make the first attempt to pre-train a 3D network with CLIP’s knowledge for 3D scene understanding.

Cross-modal Knowledge Distillation. Recently, an in-

creasing number of researchers have focused on transferring the knowledge in 2D images to 3D point cloud [34, 40]. PPKT [34] proposes the contrastive pixel-to-point knowledge transfer to utilize the rich information in image backbones. SLiDR [40] resorts to the InfoNCE loss to help the 3D network distil rich knowledge from the 2D image backbone. Our work explores leveraging the image-text pre-trained CLIP knowledge to help 3D scene understanding.

3. Methodology

Considering the impressive open-vocabulary performance achieved by CLIP in image classification and segmentation, natural curiosities have been raised. Can CLIP endow the ability to a 3D network for annotation-free scene understanding? And further, will it promote the network performance when fine-tuned on labelled data? To answer the above questions, we study the cross-modal knowledge transfer of CLIP for 3D scene understanding, termed **CLIP2Scene**. Our work is a pioneer in exploiting CLIP knowledge for 3D scene understanding. In what follows, we revisit the CLIP applied in 2D open-vocabulary classification and semantic segmentation, then present our CLIP2Scene in detail. Our approach consists of three major components: Semantic Consistency Regularization, Semantic-Guided Spatial-Temporal Consistency Regularization, and Switchable Self-Training Strategy.

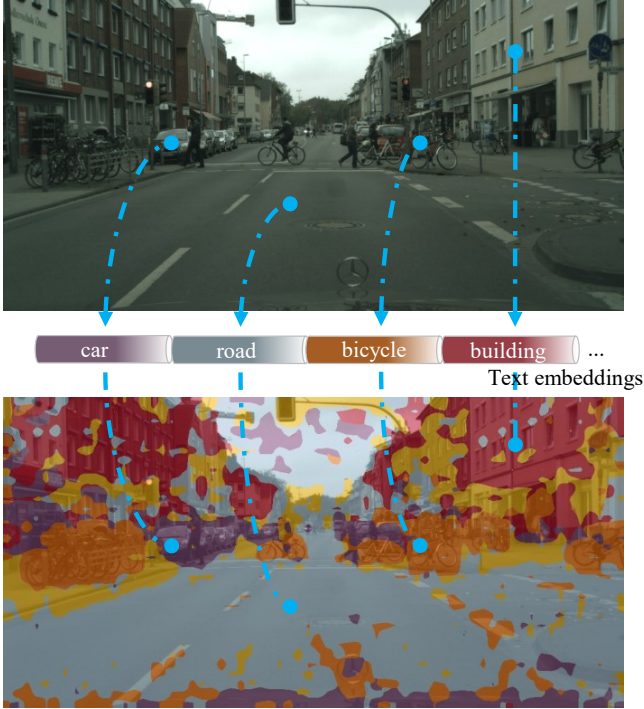


Figure 3. Illustration of the image pixel-to-text mapping. The dense pixel-text correspondence $\{x_i, t_i\}_{i=1}^M$ is extracted by the off-the-shelf method MaskCLIP [49].

3.1. Revisiting CLIP

Contrastive Vision-Language Pre-training (CLIP) mitigates the following drawbacks that dominate the computer vision field: 1. Deep models need a large amount of formatted and labelled training data, which is expensive to acquire; 2. The model’s generalization ability is weak, making it difficult to migrate to a new scenario with unseen objects. CLIP consists of an image encoder (ResNet [28] or ViT [6]) and a text encoder (Transformer [42]), both respectively project the image and text representation to a joint embedding space. During training, CLIP constructs positive and negative samples from 400 million image-text pairs to train both encoders with a contrastive loss, where the large-scale image-text pairs are free-available from the Internet and assumed to contain every class of images and most concepts of text. Therefore, CLIP can achieve promising open-vocabulary recognition.

For 2D zero-shot classification, CLIP first places the class name into a pre-defined template to generate the text embeddings and then encodes images to obtain image embeddings. Next, it calculates the similarity matrices between images and text embeddings to determine the class. MaskCLIP further extends CLIP into 2D semantic segmentation. Specifically, MaskCLIP modifies the attention pooling layer of the CLIP’s image encoder, thus performing pixel-level mask prediction instead of the global image-

level prediction.

3.2. CLIP2Scene

As shown in Fig. 2, we first leverage CLIP and 3D network to respectively extract the text embeddings, image pixel feature and point feature. Secondly, we construct positive and negative samples based on CLIP’s knowledge. Lastly, we impose Semantic Consistency Regularization by pulling the point features to their corresponding text embeddings. At the same time, we apply Spatial-Temporal Consistency Regularization by forcing the consistency between temporally coherent point features and their corresponding pixel features. In what follows, we present the details and insights.

3.2.1 Semantic Consistency Regularization

As CLIP is pre-trained on 2D images and text, our first concern is the domain gap between 2D images and the 3D point cloud. To this end, we build dense pixel-point correspondence and transfer image knowledge to the 3D point cloud via the pixel-point pairs. Specifically, we calibrate the LiDAR point cloud with corresponding images captured by six cameras. Therefore, the dense pixel-point correspondence $\{x_i, p_i\}_{i=1}^M$ can be obtained accordingly, where x_i and p_i indicates i -th paired image feature and point feature, which are respectively extracted by the CLIP’s image encoder and the 3D network. M is the number of pairs. Note that it is an online operation and is irreverent to the image and point data augmentation.

Previous methods [40, 34] provide a promising solution to cross-modal knowledge transfer. They first construct positive pixel-point pairs $\{x_i, p_i\}_{i=1}^M$ and negative pairs $\{x_i, p_j\} (i \neq j)$, and then pull in the positive pairs while pushing away the negative pairs in the embedding space via the InfoNCE loss. Despite the encourageable performance of previous methods in transferring cross-modal knowledge, they are both confronted with the same optimization-conflict issue. For example, suppose i -th pixel x_i and j -th point p_j are in the different positions of the same instance with the same semantics. However, the InfoNCE loss will try to push them away, which is unreasonable and hammer the performance of the downstream tasks [40]. In light of this, we propose a Semantic Consistency Regularization that leverages the CLIP’s semantic information to alleviate this issue. Specifically, we generate the dense pixel-text pairs $\{x_i, t_i\}_{i=1}^M$ by following the off-the-shelf method MaskCLIP [49] (Fig. 3), where t_i is the text embedding generated from the CLIP’s text encoder. Note that the pixel-text mappings are free-available from CLIP without any additional training. We then transfer pixel-text pairs to point-text pairs $\{p_i, t_i\}_{i=1}^M$ and utilize the text semantics to select the positive and negative point samples for contrastive

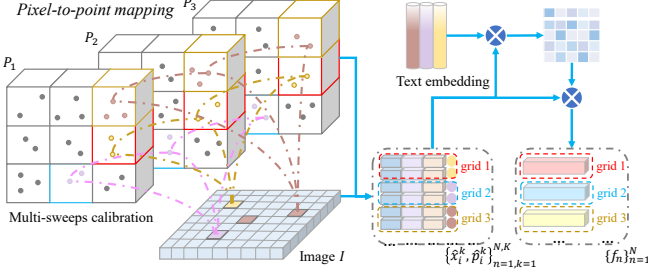


Figure 4. Illustration of the image pixel-to-point mapping (left) and semantic-guided fusion feature generation (right). We build the grid-wise correspondence between an image I and the temporally coherent LiDAR point cloud $\{P_k\}_{k=1}^K$ within S seconds and generate semantic-guided fusion features for individual grids. Both $\{\hat{x}_i^k, \hat{p}_i^k\}_{i=1, k=1}^{\hat{M}, K}$ and $\{f_n\}_{n=1}^N$ are used to perform Spatial-Temporal Consistency Regularization.

learning. The objective function is as follows:

$$\mathcal{L}_{S.info} = - \sum_{c=1}^C \log \frac{\sum_{t_i \in c, p_i} \exp(D(t_i, p_i)/\tau)}{\sum_{t_i \in c, t_j \notin c, p_j} \exp(D(t_i, p_j)/\tau)}, \quad (1)$$

where $t_i \in c$ indicates that t_i is generated by c -th classes name, and C is the number of classes. D denotes the scalar product operation and τ is a temperature term ($\tau > 0$).

Since the text is composed of class names placed into pre-defined templates, the text embedding represents the semantic information of the corresponding class. Therefore, those points with the same semantics will be restricted near the same text embedding, and those with different semantics will be pushed away. To this end, our Semantic Consistency Regularization causes less conflict in contrastive learning.

3.2.2 Semantic-guided Spatial-temporal Consistency Regularization

Besides semantic consistency regularization, we consider how image pixel features help to regularize a 3D network. The natural alternative directly pulls in the point feature with its corresponding pixel in the embedding space. However, after trial and error, we observe that the network easily degenerates and achieves poor performance in the downstream tasks when following the aforementioned strategy. The main reason lies in the noise-assigned semantics of the image pixel and the imperfect pixel-point mapping caused by the calibration errors. To this end, we propose a novel semantic-guided Spatial-Temporal Consistency Regularization to alleviate the problem by imposing a soft constraint on points within local space and time.

Specifically, given an image I and temporally coherent LiDAR point cloud $\{P_k\}_{k=1}^K$, where K is the number of sweeps within S seconds. Note that the image is matched to the first frame of the point cloud P_1 with pixel-point pairs

$\{\hat{x}_i^1, \hat{p}_i^1\}_{i=1}^{\hat{M}}$. We register the rest of the point cloud to the first frame via the calibration matrices and map them to the image (Fig. 4). Thus we obtain all pixel-point-text pairs in a short temporal $\{\hat{x}_i^k, \hat{p}_i^k, t_i^k\}_{i=1, k=1}^{\hat{M}, K}$. Next, we divide the entire stitched point cloud into regular grids $\{g_n\}_{n=1}^N$, where the temporally coherent points are located in the same grid. We impose the spatial-temporal consistency constraint within individual grids by the following objective function:

$$\mathcal{L}_{SSR} = \sum_{g_n} \sum_{(\hat{i}, \hat{k}) \in g_n} (1 - \text{sigmoid}(D(\hat{p}_i^{\hat{k}}, f_n))) / N, \quad (2)$$

where $(\hat{i}, \hat{k}) \in g_n$ indicates the pixel-point pair $\{\hat{x}_i^{\hat{k}}, \hat{p}_i^{\hat{k}}\}$ is located in the n -th grid. $\{f_n\}_{n=1}^N$ is a semantic-guided cross-modal fusion feature formulated by:

$$f_n = \sum_{(\hat{i}, \hat{k}) \in g_n} a_i^{\hat{k}} * \hat{x}_i^{\hat{k}} + b_i^{\hat{k}} * \hat{p}_i^{\hat{k}}, \quad (3)$$

where $a_i^{\hat{k}}$ and $b_i^{\hat{k}}$ are attention weight calculated by:

$$a_i^{\hat{k}} = \frac{\exp(D(\hat{x}_i^{\hat{k}}, t_i^1)/\lambda)}{\sum_{(\hat{i}, \hat{k}) \in g_n} \exp(D(\hat{x}_i^{\hat{k}}, t_i^1)/\lambda) + \exp(D(\hat{p}_i^{\hat{k}}, t_i^1)/\lambda)},$$

$$b_i^{\hat{k}} = \frac{\exp(D(\hat{p}_i^{\hat{k}}, t_i^1)/\lambda)}{\sum_{(\hat{i}, \hat{k}) \in g_n} \exp(D(\hat{x}_i^{\hat{k}}, t_i^1)/\lambda) + \exp(D(\hat{p}_i^{\hat{k}}, t_i^1)/\lambda)}, \quad (4)$$

where λ is the temperature term.

Actually, those pixel and point features within the local grid g_n are restricted near a dynamic centre f_n . Thus, such a soft constraint alleviates the noisy prediction and calibration error issues. At the same time, it imposes Spatio-Temporal Regularization on the temporally coherent point features.

3.2.3 Switchable Self-training Strategy

We combine the loss function $\mathcal{L}_{S.info}$ and \mathcal{L}_{SSR} to end-to-end train the whole network, where the CLIP's image and text encoder backbone are frozen during training. We find that method worked only when the pixel-point feature $\{x_i, p_i\}_{i=1}^M$ and $\{\hat{x}_i^k, \hat{p}_i^k\}_{i=1, k=1}^{\hat{M}, K}$, which are used in $\mathcal{L}_{S.info}$ and \mathcal{L}_{SSR} , are generated from different learnable linear layer. On top of that, we further put forward an effective strategy to promote performance. Specifically, after contrastive learning of the 3D network for a few epochs, we randomly switch the point labels between the paired image pixel's labels and their own predictions for self-training. Merely training the 3D network with their own predictions yields satisfactory performance. Essentially, such a Switchable Self-Training Strategy (S3) increases the number of

Table 1. Ablation study experiments on the nuScenes validation dataset for annotation-free semantic segmentation.

Ablation target	Settings	mIoU(%)
-	baseline	15.1
	nuScenes	15.1
Prompts	semanticKITTI	13.9
	Cityscapes	11.3
Regularization	w/o SCR	19.8
	KL	0
Training Strategies	w/o S3	18.8
	ST	10.1
	1 sweep	18.7
	3 sweeps	20.8
Sweeps	5 sweeps	20.6
	merged	18.6
-	CLIP2Scene	20.8

positive and negative samples by switching the point pseudo labels, which benefits cross-modal knowledge distillation.

4. Experiments

Datasets. We conduct experiments on two large-scale outdoor LiDAR segmentation benchmarks, *i.e.*, SemanticKITTI [3] and nuScenes [5, 22]. The nuScenes dataset contains 700 scenes for training, 150 scenes for validation and 150 scenes for testing, where 16 classes are utilized for LiDAR semantic segmentation. As to SemanticKITTI, it contains 19 classes for training and evaluation. It has 22 sequences, where sequences 00 to 10, 08 and 11 to 21 are used for training, validation and testing, respectively.

Implementation Details. We use the nuScenes [5, 22] dataset to pre-train the network. Following SLiDR, we pre-train the network on all key frames from 600 scenes. Besides, we fine-tune the pre-trained network on SemanticKITTI [3] to verify the generalization ability. We leverage CLIP’s image encoder and text encoder to generate image features and text embedding, respectively. Following MaskCLIP, we modify the attention pooling layer of the CLIP’s image encoder, thus extracting the dense pixel-text correspondences. We take SPVCNN [41] as the 3D network to produce the point-wise feature. The whole network is trained on the PyTorch platform. The training time is about 40 hours for 20 epochs on two NVIDIA Tesla A100 GPUs. For the switchable self-training strategy, we randomly switch the point supervision signal after 10 epochs. The optimizer is SGD with a cosine scheduler. We set the temperature λ and τ to be 1 and 0.5, respectively. The sweep number is set to be 3 empirically. We apply several data augmentations in contrastive learning, including random rotation around the z-axis and random flip on the

Table 2. Comparison of different self-supervised methods for semantic segmentation on the nuScenes and SemanticKITTI validation datasets.

Initialization	nuScenes		semanticKITTI
	1%	100%	1%
Random	42.2	69.1	32.5
PPKT [34]	48.0	70.1	39.1
SLiDR [40]	48.2	70.4	39.6
CLIP2Scene	56.3	71.5	42.6

point cloud, random horizontal flip and random crop-resize on the image.

4.1. Annotation-free Semantic Segmentation

After pre-training the network, we show the performance of the 3D network when it is not fine-tuned on any annotations. As no previous method reports the 3D annotation-free segmentation performance, we compare our method with different setups (Table 1). In what follows, we describe the experimental settings and give insights into our method and the different settings.

Settings. We conduct experiments on the nuScenes dataset to evaluate the annotation-free semantic segmentation performance. Following MaskCLIP [49], we place the class name into 85 hand-craft prompts and feed it into the CLIP’s text encoder to produce multiple text features. We then average the text features and feed the averaged features to the classifier for point-wise prediction. Besides, to explore how to effectively transfer CLIP’s knowledge to the 3D network for annotation-free segmentation, We conduct the following experiments to highlight the effectiveness of different modules in our framework.

Baseline. The input of the 3D network is only one sweep, and we pre-train the framework via semantic consistency regularization.

Prompts (nuScenes, semanticKITTI, Cityscapes). Based on the baseline, we respectively replace the nuScenes, semanticKITTI, and Cityscapes class names into the prompts to produce the text embedding.

Regularization (w/o STR, KL). Based on the full method, we remove the Spatial-temporal Consistency Regularization (w/o SCR). Besides, we abuse both SR and SCR and distill the image feature to the point cloud by Kullback–Leibler (KL) divergence loss.

Training Strategies (w/o S3, ST). We abuse the Switchable Self-Training Strategy (w/o S3) in the full method. Besides, we show the performance of only training the 3D network by their own predictions after ten epochs (ST).

Sweeps Number (1 sweep, 3 sweeps, 5 sweeps, and merged). We set the sweep number K to be 1, 3, and 5, respectively. Besides, we also take three sweeps of the point cloud as the input to pre-train the network.

Effect of Different Prompts. To verify how text em-

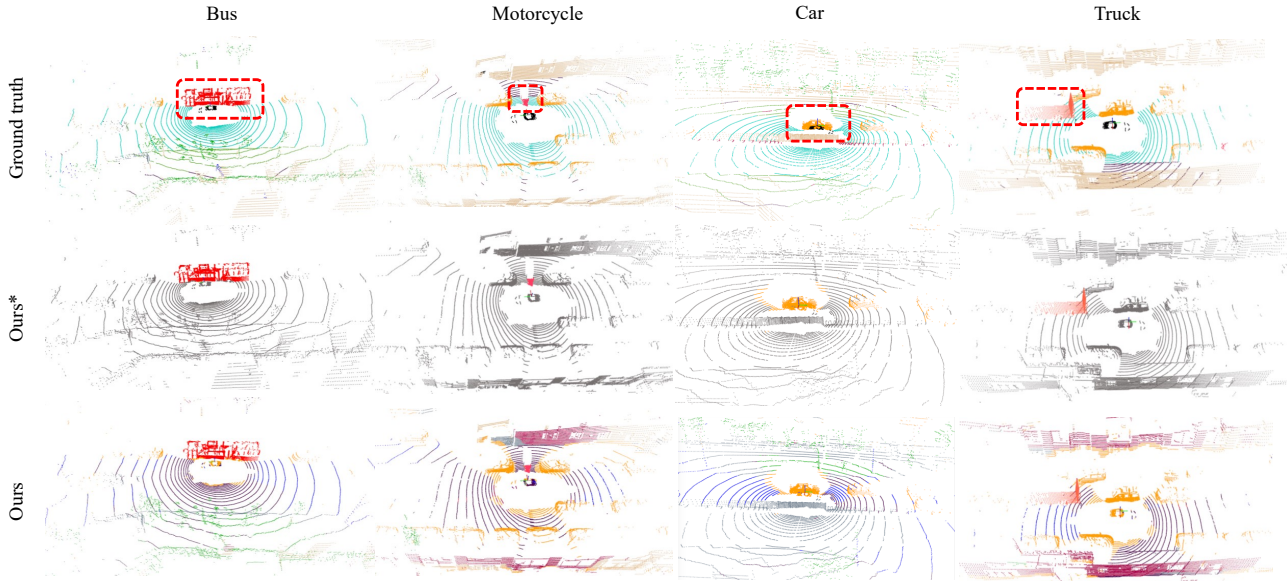


Figure 5. Qualitative results of annotation-free semantic segmentation on nuScenes dataset. Note that we show the results by individual class. From the left to the right column are the bus, motorcycle, car and truck, respectively. The first row is the ground truth; The second row (ours*) is our prediction of the highlighted target; the third row is our prediction of full classes (ours).

bedding affects the performance, we generate various text embeddings by the class name from different datasets (nuScenes, SemanticKIT, and Cityscapes) for pre-training the framework. As shown in Table 1, we find that even learning with other datasets’ text embedding (semanticKIT and Cityscapes), the 3D network could still recognize the nuScenes’s objects with decent performance (13.9 and 11.3 mIoU, respectively). The result shows that the 3D network is capable of open-vocabulary recognition.

Effect of Semantic and Spatial-temporal Consistency Regularization. We remove Spatial-temporal Consistency Regularization (w/o SCR) from our method. Experiments show that the performance is dramatically decreased, indicating the effectiveness of our design. Besides, we also distill the image feature to the point cloud by KL divergence loss, where the text embeddings calculate the logits. However, such a method fails to transfer the semantic information from the image. The main reason is the noise-assigned semantics of the image pixel and the imperfect pixel-point correspondence due to the calibration error.

Effect of Switchable Self-training Strategy. To examine the effect of the Switchable Self-Training Strategy, we either train the network with image supervision (w/o S3) or train the 3D network by their own predictions. Both trials witness the performance drop, indicating our Switchable Self-Training Strategy is efficient in cross-modal self-supervised learning. The main reason is that the number of positive and negative samples is enlarged by switching the supervision signal.

Effect of Sweep Numbers. Intuitively, the performance of our method benefits from more sweeps information. Therefore, we also show the performance when restricting sweep size to 1, 3, and 5, respectively. However, we observe that the performance of 5 sweeps is similar to 3 sweeps but is more computationally expensive. Thus, we empirically set the sweep number to be 3.

Qualitative Evaluation. We show the qualitative evaluation in Fig. 5. Note that we show the results by individual class (construction vehicle, truck, and car). The results show that our method is able to perceive the objects without any annotation training data. However, we also observe the false positive predictions around the ground truth objects. We will resolve this issue in future work.

4.2. Annotation-efficient Semantic Segmentation

Besides annotation-free semantic segmentation, the pre-trained 3D network also boosts the performance when it is fine-tuned on labelled data. To the best of our knowledge, only one published method SLidR studies image-to-Lidar self-supervised representation distillation. We also compared our method with another self-supervised method PPKT [34] for 3D network pre-training. In the followings, we first introduce SLidR [40] and PPKT, then compare them in detail.

PPKT. PPKT is a cross-modal self-supervised method for the RGB-D dataset. It performs 2D-to-3D knowledge distillation via pixel-to-point contrastive loss. Since there is no public code, we re-implement it for a fair comparison.

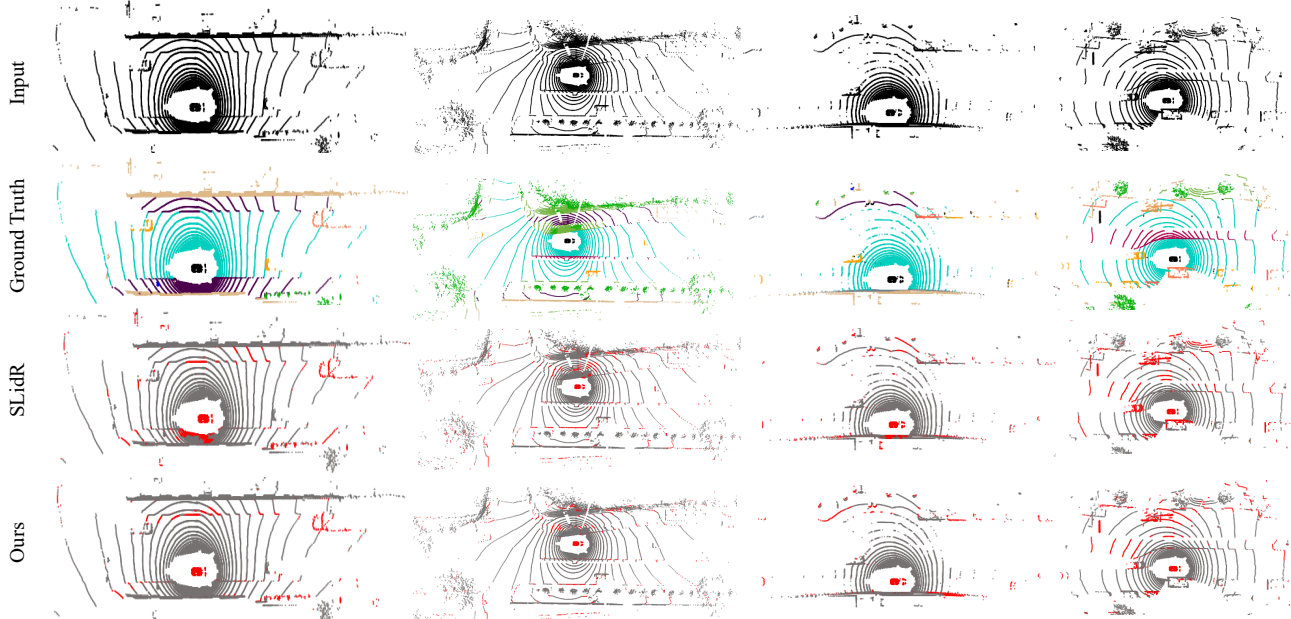


Figure 6. Qualitative results of fine-tuning on 1% nuScenes dataset. From the first row to the last row are the input Lidar scan, ground truth, prediction of SLidR, and our prediction, respectively. Note that we show the results by error map, where the red point indicates the wrong prediction. Apparently, our method achieves decent performance.

Specifically, we use the same 3D network and training protocol but replace our semantic and Spatio-Temporal Regularization with InfoNCE loss. The framework is trained on 4, 096 randomly selected image-to-point pairs for 50 epochs.

SLidR. SLidR is an image-to-Lidar self-supervised method for autonomous driving data. Compared with PPKT, it introduces image super-pixel into cross-modal self-supervised learning. For a fair comparison, we replace our loss function with their superpixel-driven contrastive loss.

Performance. As shown in Table 2, our method significantly outperforms the state-of-the-art methods when fine-tuned on 1% and 100% data, with the improvement of 8.1% and 1.1%, respectively. Compared with the random initialization, the improvement is 14.1% and 2.4%, respectively, indicating the efficiency of our semantic-driven cross-modal contrastive learning framework. The qualitative results are shown in Fig. 6. Besides, we also verify the cross-domain generalization ability of our method. When pre-training the 3D network on the nuScenes dataset and fine-tuning on 1% SemanticKITTI dataset, our method significantly outperforms other state-of-the-art self-supervised methods.

Discussions. PPKT and SLidR reveal that contrastive loss is promising for transferring knowledge from image to point cloud. Like self-supervised learning, constructing the positive and negative samples is vital to unsupervised cross-modal knowledge distillation. However, previous methods suffer from the optimization-conflict issue, i.e., some

of the negative paired samples are actually positive pairs. For example, the road occupies a large proportion of the point cloud in a scene and is supposed to have the same semantics in the semantic segmentation task. When randomly selecting training samples, most negatively defined road-road points are actually positive. When feedforwarding such training samples into contrastive learning, the contrastive loss will push them away in the embedding space, leading to unsatisfactory representation learning and hampering the downstream tasks’ performance. SLidR introduces superpixel-driven contrastive learning to alleviate such issues. The motivation is that the visual representation of the image pixel and the projected points are consistent intra-superpixel. Although avoiding selecting the negative image-point pairs from the same superpixel, the conflict issue still exists inter-superpixel. In our CLIP2Scene, we introduce the free-available dense pixel-text correspondence to alleviate the optimization conflicts. The text embedding represents the semantic information and can be used to select more reasonable training samples for contrastive learning.

Besides training sample selection, the previous method also ignores the temporal coherence of the multi-sweep point cloud. Similar to multi-view consistency, multi-sweep consistency emphasizes inter-sweep consistency along time series. That is, for those LiDAR points mapping to the same image pixel, their feature should be the same. Besides, considering the sparsity of the LiDAR scan and the calibration error between the LiDAR scan and the camera image. We

relax the pixel-to-point mapping to image grid-to-point grid mapping and calculate the dynamic centre within the individual grid for consistency regularization. To this end, our Spatial-temporal consistency regularization leads to a more comprehensive point representation.

Last but not least, the previous method typically enlarges the number of training samples by data augmentation. In our CLIP2Scene, we find that randomly switching the supervision signal benefits self-supervised learning. Essentially, our Switchable Self-Training Strategy enlarges the training samples and prevents the network from deteriorating.

5. Conclusion

We explored how CLIP knowledge benefits 3D scene understanding in this paper, termed CLIP2Scene. To efficiently transfer CLIP’s image feature and text feature to a 3D network, we propose a novel Semantic-driven Cross-modal Contrastive Learning framework including Semantic Regularization and Spatial-Temporal Regularization. For the first time, our pre-trained 3D network achieves annotation-free 3D semantic segmentation with decent performance. Besides, our method significantly outperforms state-of-the-art self-supervised methods when fine-tuning the 3D network with labelled data.

Potential Negative Impacts. Although our approach improves the 3D semantic segmentation performance in general, its effectiveness under adversarial attack is not considered, which could be safety-critical in practical applications, such as autonomous driving and robot navigation.

References

- [1] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid. Label-embedding for attribute-based classification. *CVPR*, pages 819–826, 2013. 2
- [2] Z. Akata, S. E. Reed, D. Walter, H. Lee, and B. Schiele. Evaluation of output embeddings for fine-grained image classification. *CVPR*, pages 2927–2936, 2015. 2
- [3] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall. SemanticKITTI: A Dataset for Semantic Scene Understanding of Lidar Sequences. In *IEEE International Conference on Computer Vision*, pages 9297–9307, 2019. 6
- [4] M. Bucher, S. Herbin, and F. Jurie. Generating visual representations for zero-shot classification. *ICCVW*, pages 2666–2673, 2017. 2
- [5] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom. NuScenes: A Multimodal Dataset for Autonomous Driving. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 11621–11631, 2020. 6
- [6] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 4
- [7] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020. 3
- [8] S. Changpinyo, W.-L. Chao, B. Gong, and F. Sha. Synthesized classifiers for zero-shot learning. *CVPR*, pages 5327–5336, 2016. 2
- [9] N. Chen, L. Liu, Z. Cui, R. Chen, D. Ceylan, C. Tu, and W. Wang. Unsupervised learning of intrinsic structural representation points. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9121–9130, 2020. 3
- [10] R. Chen, Z. Xinge, N. Chen, D. Wang, W. Li, Y. Ma, R. Yang, and W. Wang. Referring self-supervised learning on 3d point cloud. 2021. 3
- [11] R. Chen, X. Zhu, N. Chen, W. Li, Y. Ma, R. Yang, and W. Wang. Zero-shot point cloud segmentation by transferring geometric primitives. *arXiv preprint arXiv:2210.09923*, 2022. 2
- [12] R. Chen, X. Zhu, N. Chen, D. Wang, W. Li, Y. Ma, R. Yang, and W. Wang. Towards 3d scene understanding by referring synthetic models. *arXiv preprint arXiv:2203.10546*, 2022. 3
- [13] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 3
- [14] X. Chen and K. He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021. 3
- [15] R. Cheng, R. Razani, E. Taghavi, E. Li, and B. Liu. (af)2-S3Net: Attentive Feature Fusion with Adaptive Feature Selection for Sparse Semantic Segmentation Network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 12547–12556, 2021. 2
- [16] A. Cheraghian, S. Rahman, D. Campbell, and L. Petersson. Mitigating the Hubness Problem for Zero-Shot Learning of 3D Objects. *arXiv preprint arXiv:1907.06371*, 2019. 2
- [17] A. Cheraghian, S. Rahman, T. F. Chowdhury, D. Campbell, and L. Petersson. Zero-Shot Learning on 3D Point Cloud Objects and Beyond. *International Journal of Computer Vision*, 130(10):2364–2384, 2022. 2
- [18] A. Cheraghian, S. Rahman, and L. Petersson. Zero-Shot Learning of 3D Point Cloud Objects. In *International Conference on Machine Vision Applications (MVA)*, pages 1–6. IEEE, 2019. 2
- [19] B. Demirel, R. G. Cinbis, and N. Ikizler-Cinbis. Attributes2classname: A discriminative model for attribute-based unsupervised zero-shot learning. *ICCV*, pages 1241–1250, 2017. 2
- [20] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015. 3

- [21] A. Dosovitskiy, J. T. Springenberg, M. Riedmiller, and T. Brox. Discriminative unsupervised feature learning with convolutional neural networks. *Advances in neural information processing systems*, 27, 2014. 3
- [22] W. K. Fong, R. Mohan, J. V. Hurtado, L. Zhou, H. Caesar, O. Beijbom, and A. Valada. Panoptic nuScenes: A Large-Scale Benchmark for LiDAR Panoptic Segmentation and Tracking. *IEEE Robotics and Automation Letters*, 7:3795–3802, 2022. 6
- [23] C. Gan, M. Lin, Y. Yang, Y. Zhuang, and A. Hauptmann. Exploring semantic inter-class relationships (sir) for zero-shot action recognition. In *AAAI*, 2015. 2
- [24] B. Gao, Y. Pan, C. Li, S. Geng, and H. Zhao. Are We Hungry for 3D LiDAR Data for Semantic Segmentation? A Survey of Datasets and Methods. *IEEE Transactions on Intelligent Transportation Systems*, 2021. 1
- [25] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. 3
- [26] Y. Guo, H. Wang, Q. Hu, H. Liu, L. Liu, and M. Benamoun. Deep Learning for 3D Point Clouds: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(12):4338–4364, 2020. 1
- [27] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 3
- [28] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 4
- [29] Y. Hou, X. Zhu, Y. Ma, C. C. Loy, and Y. Li. Point-to-Voxel Knowledge Distillation for LiDAR Semantic Segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 8479–8488, 2022. 2
- [30] E. Kodirov, T. Xiang, and S. Gong. Semantic autoencoder for zero-shot learning. *CVPR*, pages 4447–4456, 2017. 2
- [31] C. H. Lampert, H. Nickisch, and S. Harmeling. Learning to detect unseen object classes by between-class attribute transfer. *CVPR*, pages 951–958, 2009. 2
- [32] C. H. Lampert, H. Nickisch, and S. Harmeling. Attribute-based classification for zero-shot visual object categorization. *IEEE TPAMI*, 36:453–465, 2014. 2
- [33] Y. Li, Z. Jia, J. Zhang, K. Huang, and T. Tan. Deep semantic structural constraints for zero-shot learning. In *AAAI*, 2018. 2
- [34] Y.-C. Liu, Y.-K. Huang, H.-Y. Chiang, H.-T. Su, Z.-Y. Liu, C.-T. Chen, C.-Y. Tseng, and W. H. Hsu. Learning From 2D: Contrastive Pixel-to-Point Knowledge Transfer for 3D Pretraining. *arXiv preprint arXiv:2104.04687*, 2021. 2, 3, 4, 6, 7
- [35] B. Michele, A. Boulch, G. Puy, M. Bucher, and R. Marlet. Generative zero-shot learning for semantic segmentation of 3d point clouds. In *3DV*, pages 992–1002. IEEE, 2021. 2
- [36] A. Mishra, M. S. K. Reddy, A. Mittal, and H. A. Murthy. A generative model for zero shot learning using conditional variational autoencoders. *CVPRW*, pages 2269–22698, 2018. 2
- [37] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 652–660, 2017. 2
- [38] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 2
- [39] Y. Rao, W. Zhao, G. Chen, Y. Tang, Z. Zhu, G. Huang, J. Zhou, and J. Lu. Denseclip: Language-guided dense prediction with context-aware prompting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18082–18091, 2022. 3
- [40] C. Sautier, G. Puy, S. Gidaris, A. Boulch, A. Bursuc, and R. Marlet. Image-to-Lidar Self-Supervised Distillation for Autonomous Driving Data. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9891–9901, 2022. 2, 3, 4, 6, 7
- [41] H. Tang, Z. Liu, S. Zhao, Y. Lin, J. Lin, H. Wang, and S. Han. Searching Efficient 3D Architectures with Sparse Point-Voxel Convolution. In *European Conference on Computer Vision*, pages 685–702. Springer, 2020. 6
- [42] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 4
- [43] Y. Xian, Z. Akata, G. Sharma, Q. N. Nguyen, M. Hein, and B. Schiele. Latent embeddings for zero-shot classification. *CVPR*, pages 69–77, 2016. 2
- [44] J. Xu, R. Zhang, J. Dou, Y. Zhu, J. Sun, and S. Pu. RPNnet: A Deep and Efficient Range-Point-Voxel Fusion Network for Lidar Point Cloud Segmentation. In *IEEE International Conference on Computer Vision*, pages 16024–16033, October 2021. 2
- [45] X. Yan, J. Gao, C. Zheng, C. Zheng, R. Zhang, S. Cui, and Z. Li. 2DPASS: 2D Priors Assisted Semantic Segmentation on LiDAR Point Clouds. In *European Conference on Computer Vision*, 2022. 2
- [46] L. Yao, J. Han, Y. Wen, X. Liang, D. Xu, W. Zhang, Z. Li, C. Xu, and H. Xu. Detclip: Dictionary-enriched visual-concept paralleled pre-training for open-world detection. *arXiv preprint arXiv:2209.09407*, 2022. 3
- [47] É. Zablocki, P. Bordes, B. Piwowarski, L. Soulier, and P. Gallinari. Context-aware zero-shot learning for object recognition. *ArXiv*, abs/1904.12638, 2019. 2
- [48] R. Zhang, Z. Guo, W. Zhang, K. Li, X. Miao, B. Cui, Y. Qiao, P. Gao, and H. Li. PointCLIP: Point Cloud Understanding by CLIP. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8552–8562, 2022. 2, 3
- [49] C. Zhou, C. C. Loy, and B. Dai. Extract Free Dense Labels from CLIP. In *European Conference on Computer Vision*, pages 696–712. Springer, 2022. 2, 3, 4, 6
- [50] X. Zhu, H. Zhou, T. Wang, F. Hong, Y. Ma, W. Li, H. Li, and D. Lin. Cylindrical and Asymmetrical 3D Convolution Networks for Lidar Segmentation. In *IEEE Conference*

on Computer Vision and Pattern Recognition, pages 9939–9948, 2021. [2](#)